

# 企业大数据技术及应用

Enterprise Big Data: Technology & Applications

塔布 / 数说故事创始人  
徐亚波博士

CONTENTS  
目录

- 01 关于塔布数说
- 02 企业大数据现状与趋势
- 03 大数据核心技术
- 04 大数据应用案例



01

# 关于塔布数说

- ◆ 塔布数说简史
- ◆ 塔布数说定位

# 塔布数说简史



## Academia

DM Models & Algo

Search & Mining

Cloud Computing

NLP & Text Mining

Data Infrastructure &

Data Integration &

Semantic Analysis &

User Modeling &

Recommendation

Big Data Services

2001

2007

2009

2010

2012

2013

2014

2015

2016

## Industry

摘要式答案引擎

新闻聚合引擎

购物搜索

微博分析工具

社会化媒体数据平台

数说故事(DATASTORY)

大数据应用

塔布(DATATUB):

大数据管理平台

→企业大数据整体服务提供商

## Data

100M~1G 结构化数据

10G~ 网页数据

1T+ 新闻数据

1T+ 商品/评论数据

10T~ Social数据

日更新1G+

企业内部+外部数据

100T ~ 1P

结构化+非结构化

日更新 1T+

日更新 100T+

# 塔布数说定位：企业大数据整体解决方案提供商



# 02

## 企业大数据现状与趋势

- ◆ 为什么企业会有大数据？
- ◆ 企业大数据有多少？
- ◆ 行业大数据在什么阶段？
- ◆ 企业大数据需求阶段

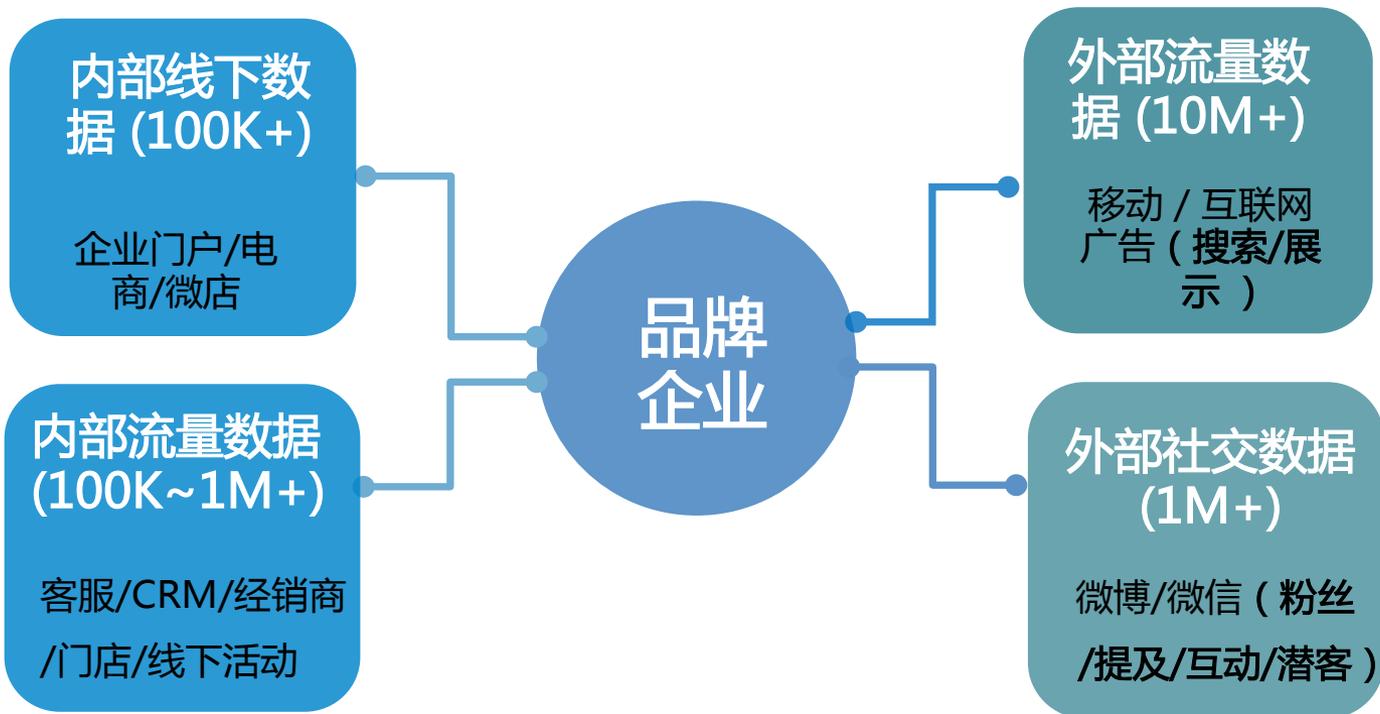
# 为什么会有企业大数据?



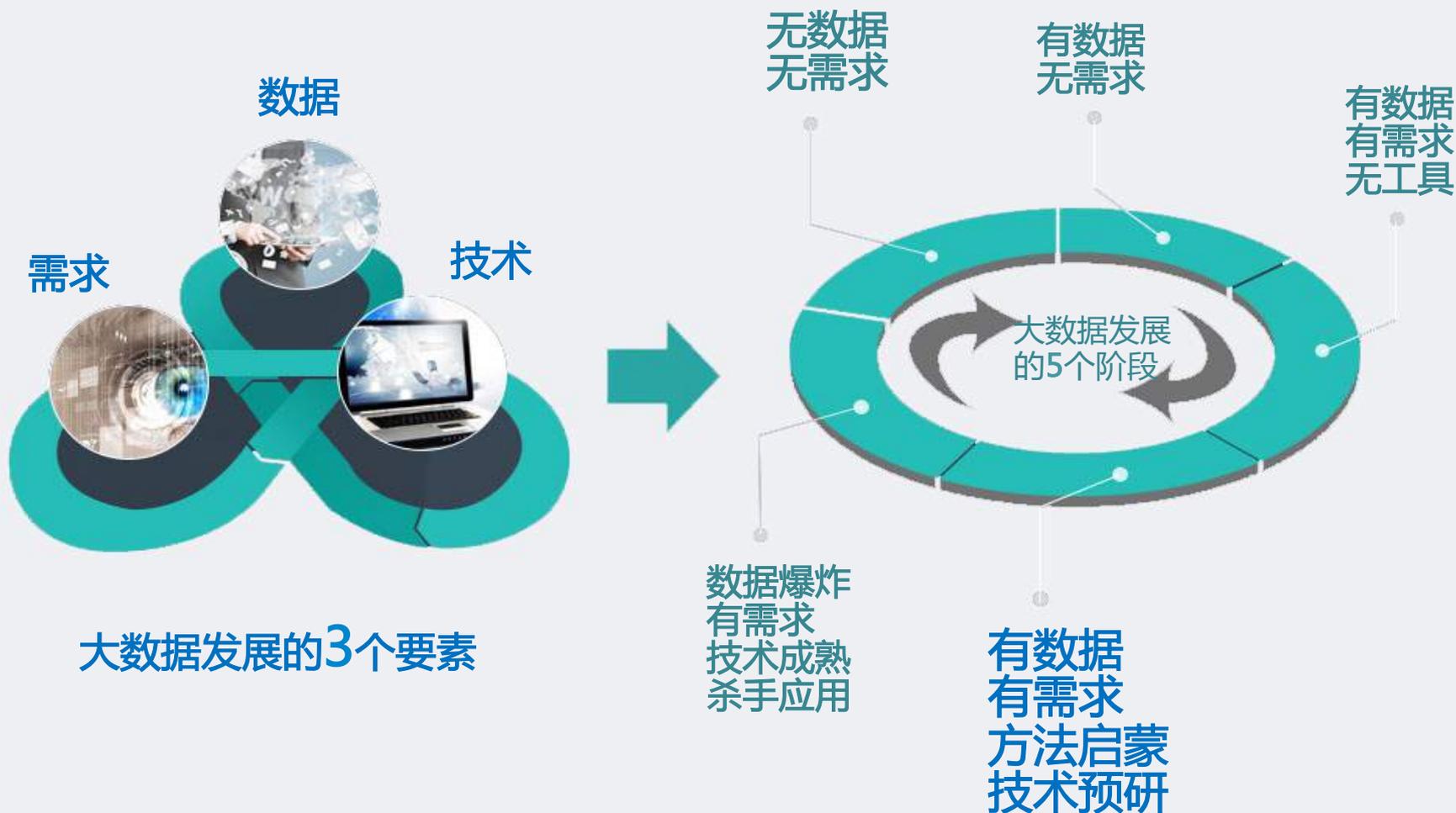
社交 / 移动 / 电商 / 数字广告

企业→消费者: 距离缩短, 接触面扩大, 核心环节数字化

# 企业有多少数据？



# 行业大数据在什么阶段？

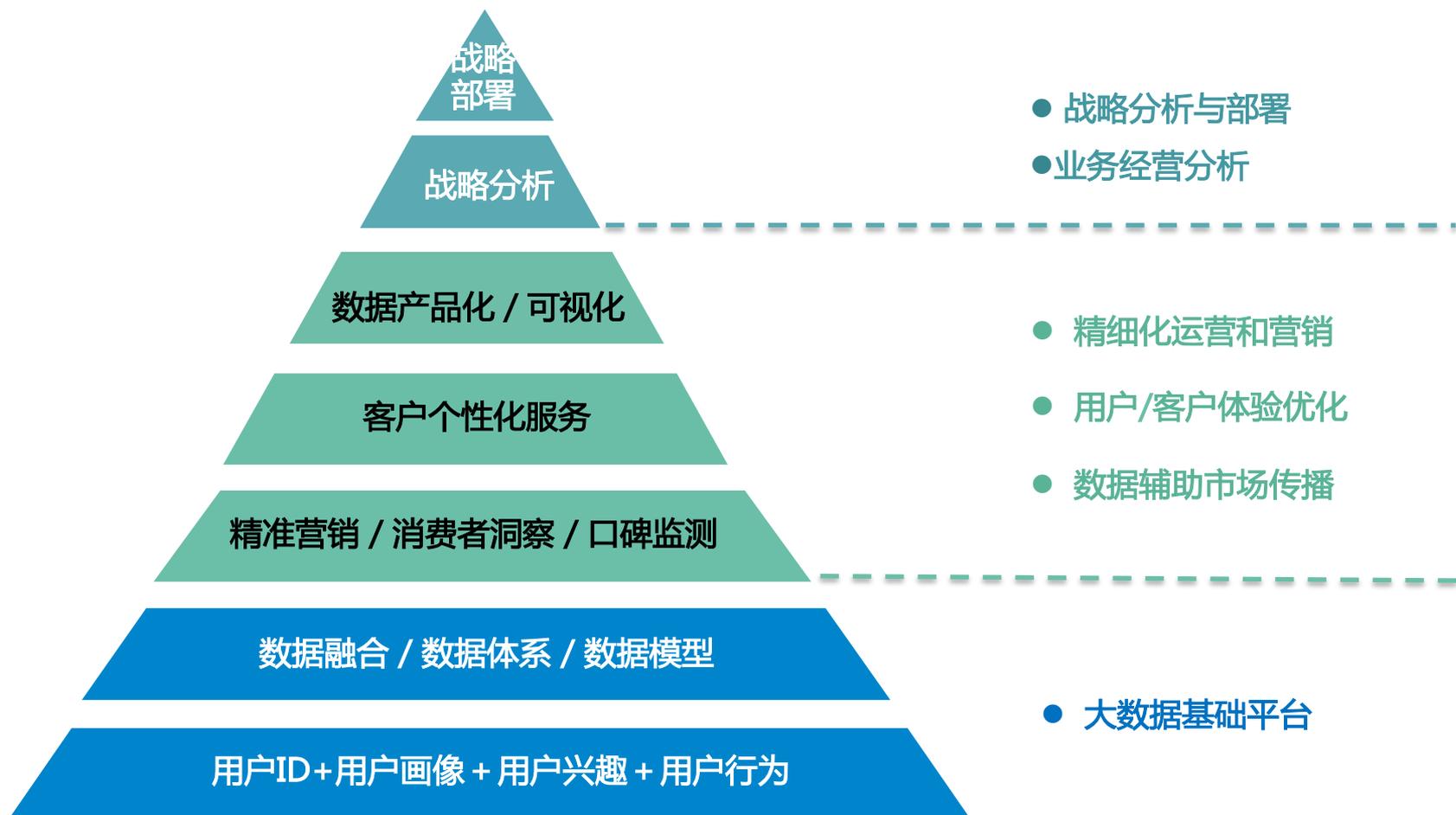


# 03

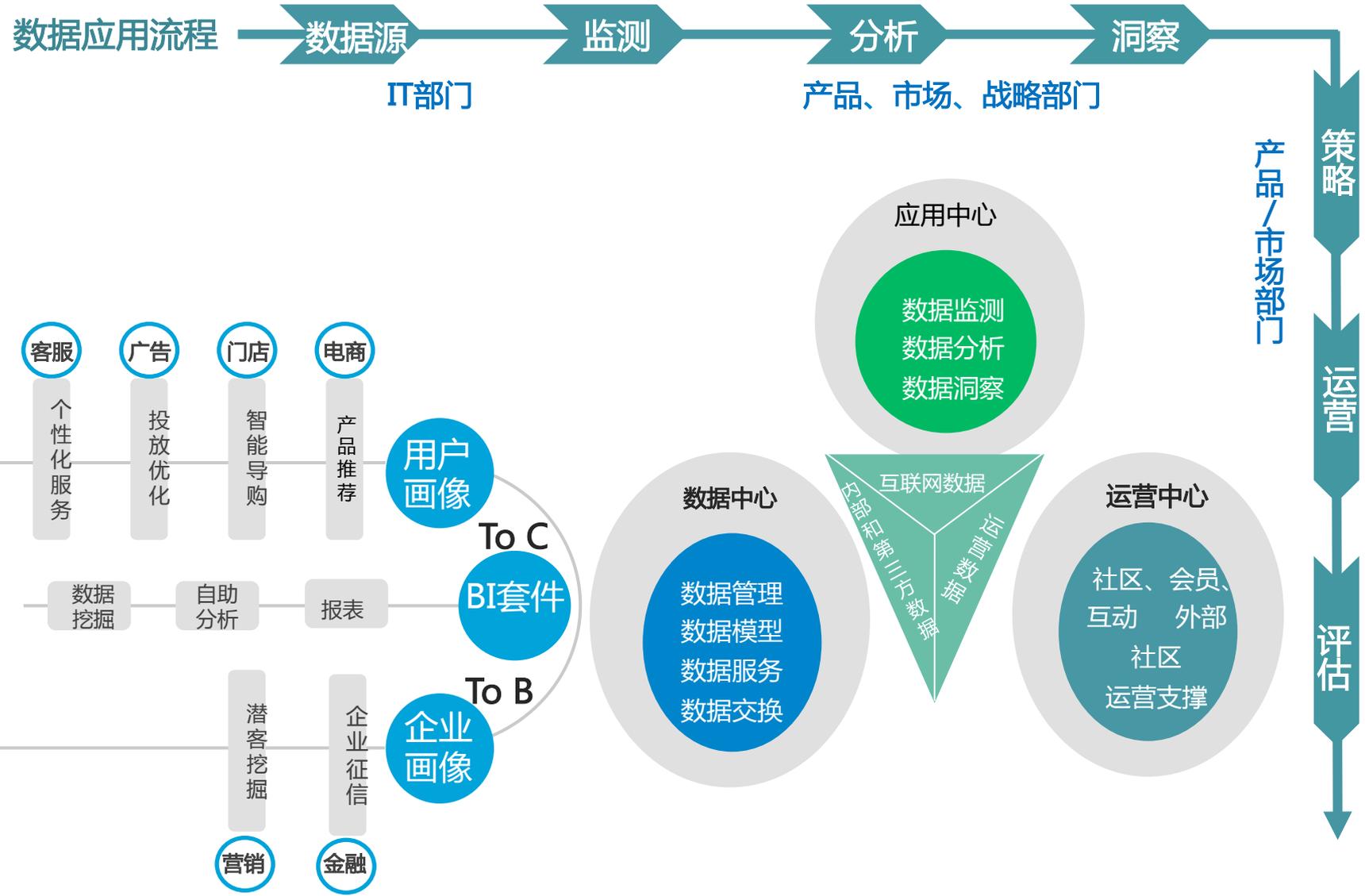
## 大数据核心技术

- ◆ 大数据应用体系
- ◆ 大数据产品体系
- ◆ 大数据平台架构
- ◆ 大数据核心技术要素

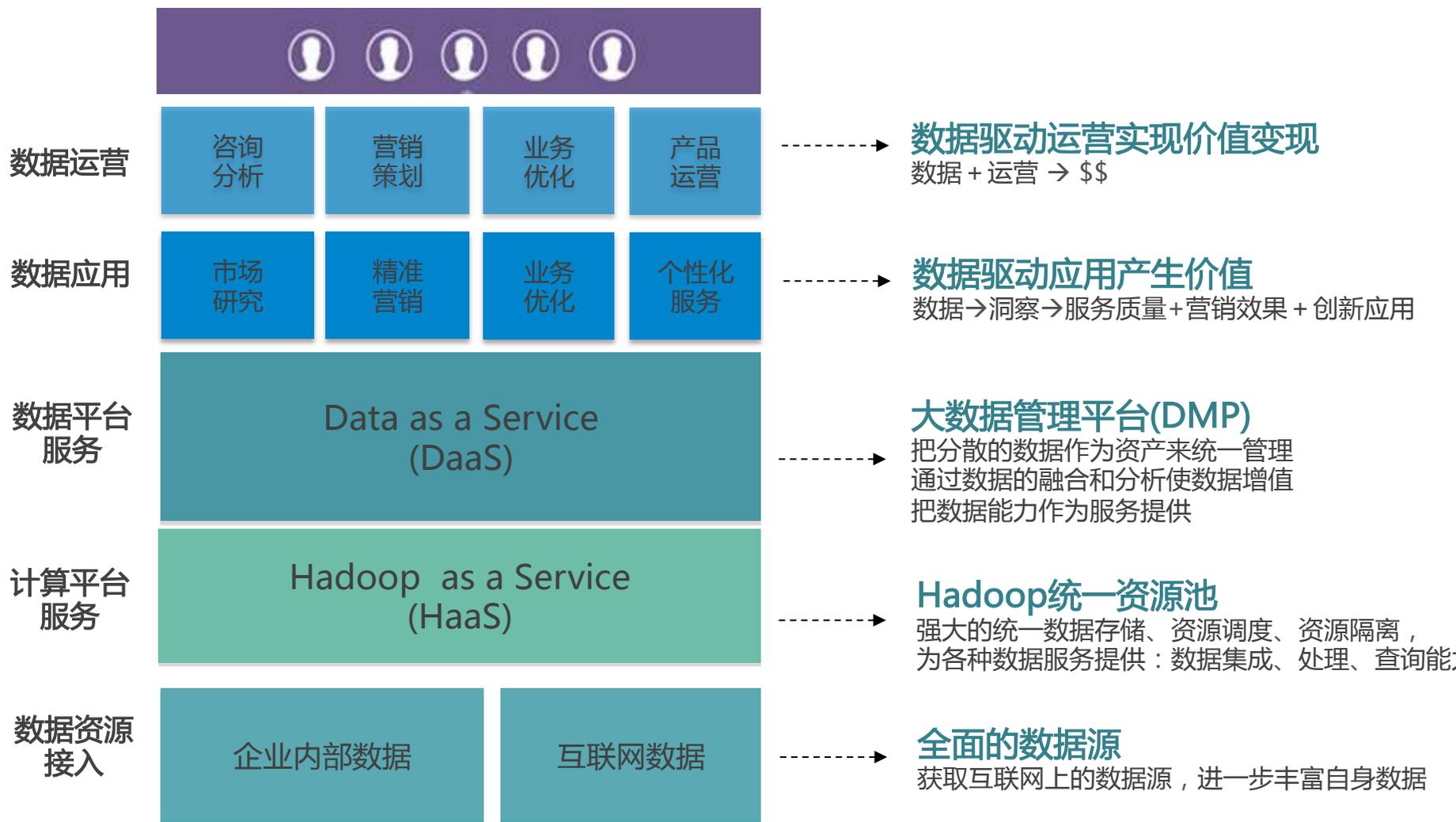
# 企业大数据应用金字塔



# 企业大数据产品体系



# 高大上的大数据平台架构



# HaaS提供统一计算和存储平台

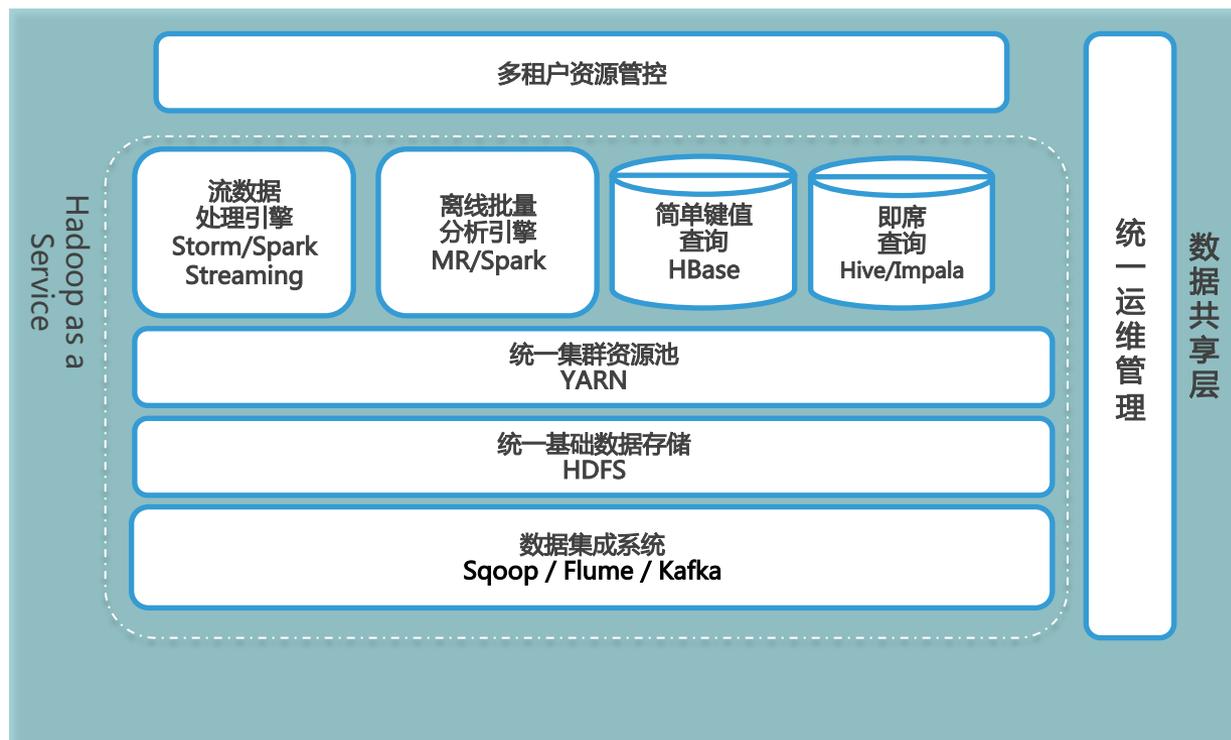
- Hadoop生态圈的优势
  - 开源的技术框架**大大降低了大数据处理的门槛和成本**，提供简洁有效的分布式计算存储的实现方式。
  - 生态圈日益完备，**渐成业内标准**

## 核心技术点：

□ **技术理解 & 规划**：Hadoop生态日趋复杂，需要业务和技术结合进行合理的架构和资源规划。

□ **统一运维管理**：开源架构、建设成本低，但技术支持不完备、更需要专业运维，同时持续跟进开源社区。

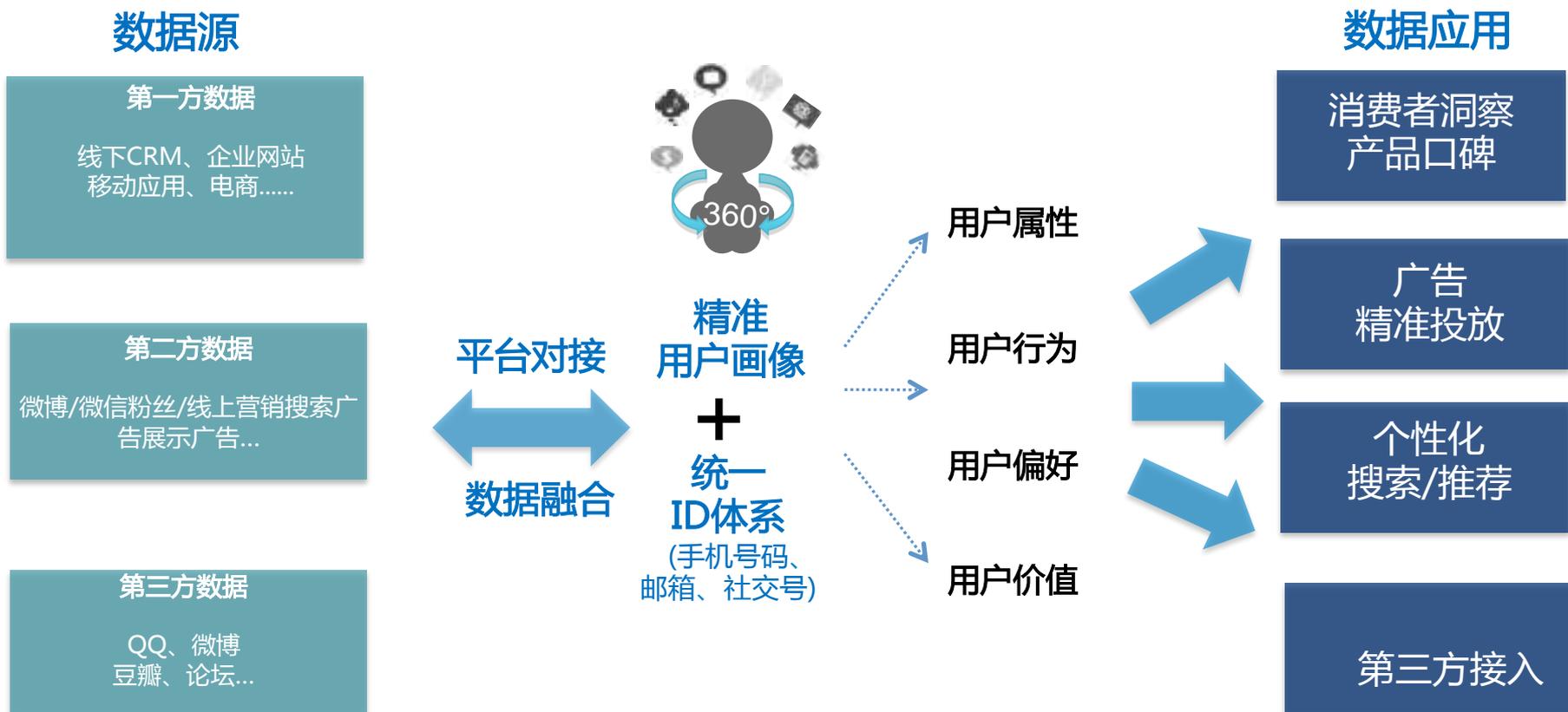
□ **多租户资源管控**：同一Hadoop集群需要支持不同业务需求，需要做好资源管控、任务优先级管理以及数据的隔离和共享



# 统一数据管理平台



## Unified Data Management Platform (UDMP)



# 数据体系规划 & 平台构建



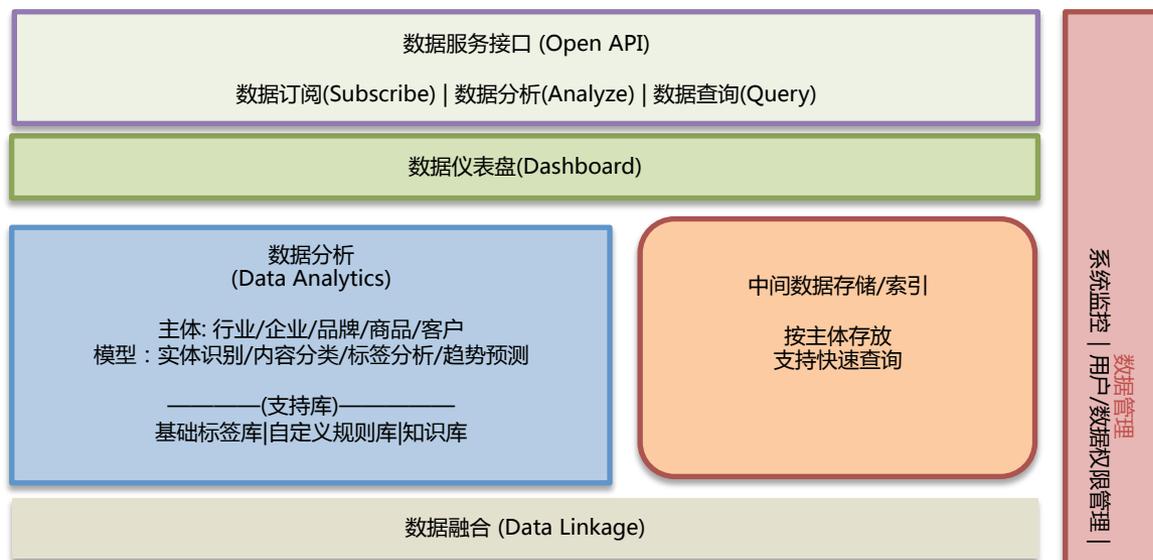
## 应用数据层

针对特定应用和问题的数据  
按应用存放



## 中间数据层

可再生数据资源，共性数据需求  
基础层之上加工的数据，根据不同业务需求，按不同主体存放  
通过数据管理平台，提供给应用使用



## 基础数据层

不可再生数据资源  
原始数据，不做汇总，应用不直接访问，仅做简单ETL和数据清洗



# 企业大数据核心技术要素



如何获得丰富的数据资源？



互联网数据采集

如何处理分散而复杂的海量数据？



海量作业平台

如何深度分析挖掘让数据具有价值？



语义分析

如何以用户为中心构建数据管理平台？



用户画像



## 核心技术点

- 多租户/多实例的使用模式; 自主开发, 基于Spark进行分布式批量/实时数据抓取
- 爬虫任务请求级别最小力度拆分, 数据DAG, 防止单个环节瓶颈
- 基于yarn的网络IO隔离, 完善的告警、监控
- 支持各种反爬策略
  - 代理IP服务, 统一调度代理IP使用
  - 移动端token获取
  - 自动化的反爬探测, 选择合理的反爬规则





## 核心技术要点

- 基于计算平台之上 JOB (任务)+TABLE(仓库) 的管理
- 丰富的内置任务类型
- 完善 workflow 调度以及失败重试
- 支持自定义任务满足企业需求
- 良好的结合 yarn 队列，对多租户进行管理
- 配套提供完善任务监控机制、使用统计



## 应用场景

多租户、复杂数据处理流程，单一框架无法满足需求，大量数据处理任务需要管理协调。

## 应用案例——

- H5数据监测: 每天亿级H5页面PV的分析处理、关系传播图绘制、KOL的发现、用户画像
- 社交媒体数据分析应用: 亿级用户数据分析，百亿级关系图，日处理百万级微博用户数据更新
- 运营商DPI数据应用: 每天百亿级移动用户访问记录的分析 and 建模

# 语义分析

## 核心技术

覆盖语义分析主要基础环节

→词级分析

词向量模型(Word2Vec)

词分类/实体识别/词义消歧/关键词抽取

→句子级分析

句子向量模型(ParagraphVec) /句法分析模型

短文本聚类/情感分析/观点抽取

→文章级分析

内容分类/内容聚类

(单文本/多文本)内容摘要

## 技术要点

长期知识库/语料的积累

基于深度学习的核心语义模型(词/句子向量模型)

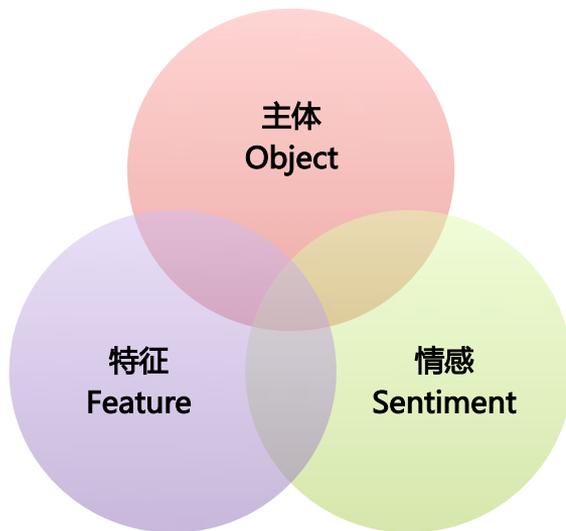
## 应用场景

各种非结构化数据分析





从任意非结构化自然语言中，提取句子最精要的观点，即OFS三元组。



今天在山姆店买披萨，服务员(Object)态度(Feature)很不好(Sentiment)

附近就这一家沃尔玛(Object)，东西(Feature)比较齐全(Sentiment)

# 传统抽取方法 vs 三元组技术



传统抽取：酒店评价的简单关键词



三元组：主体-特征-情感 的观点搭配



正面

负面

(说明：因为评价主体都是酒店，所以词云图中为“特征-情感”二元组)

传统抽取：中文歧义普遍存在

苹果 = 知名品牌？一种水果

沃尔玛服务员大大的赞  
= 服务员大？服务员-赞

声音大 = 正面评价？负面评价

三元组：解决实体、特征、情感歧义

上下文词向量模型，消除实体歧义

语法分析、相关性分类器，确定特征搭配

极性词典，解决正负面歧义

传统抽取：词汇受限于人的经验

厕所、洗手间、卫生间、男  
厕、女厕

还有什么词？

三元组：词向量模型自动补全近义词

厕所间、蹲坑、蹲位、女厕所、男厕所、  
公共厕所、公厕、洗手池、浴室、洗澡  
间、马桶、洗漱间、洗手盆、内急、抽  
水马桶、洗手台、洗澡房、小解、洗手  
盘、蹲厕所、茅房、更衣间、解手、盥  
洗室、便池、小便池、wc、如厕

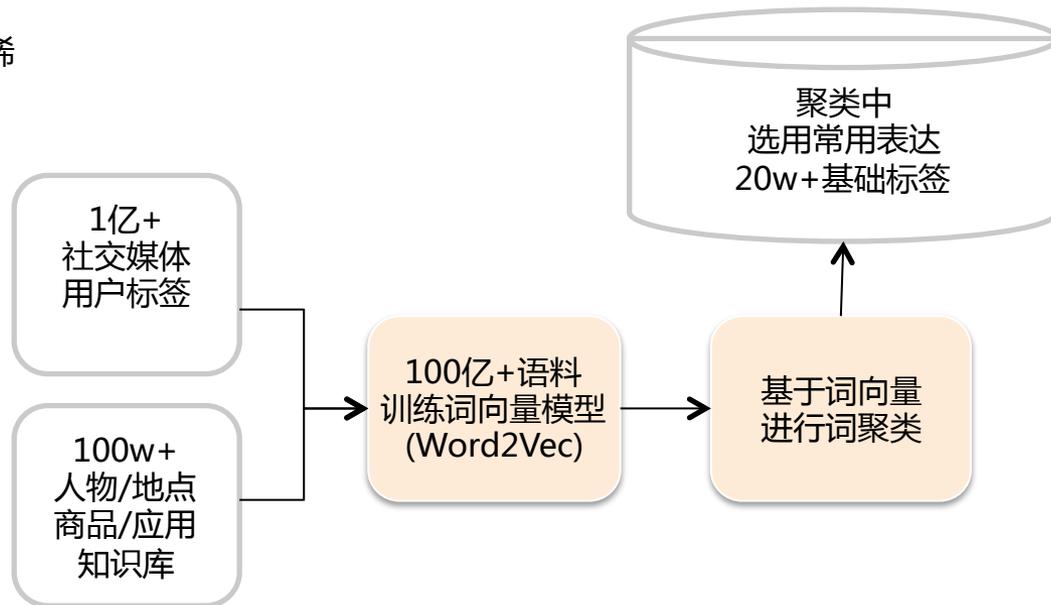
# 用户画像: 标签体系

## 传统方法问题

- 固定的标签集合很难满足业务的弹性
- 上层标签太通用, 没法描述特定细分人群
- 简单关键词模型对通用类别词效果不佳 (非常用表达)
- 高质量多层多分类模型训练数据稀缺, 类别太多人工标注也很困难

## 标签体系要求

- 要有足够的覆盖面
- 要有足够的表达能力
- 要有足够细粒度
- 要能理解其背后的语义



# 用户画像: 标签模型



## 个人属性

年龄/性别, 行业/职业, 人生状态/个性  
固定维度, 类别较少, 上层概念, 独立模型

## 行为属性

网站/APP/新闻/视频

## 位置属性

行为轨迹 + 地点

## 兴趣属性

商品/品牌/美食/影视/图书, 维度巨大, 有层次结构,  
实体识别/知识库支持

## 关系属性

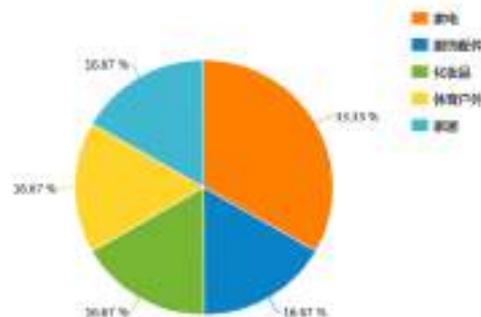
和谁交往 / 关注谁  
关系分析

## 购物属性

商品库 + 商品识别/分类

## 价值属性

购买力/影响力/活跃度  
构建特定指标及对应分类模型





04

# 大数据应用案例

# 大数据应用：数据驱动的零售渠道管理优化



利用互联网用户评论数据进行社群聆听，监控与该日化企业合作的50个零售商店相关的用户评论，**通过线上数据进行渠道/购物者研究并指导渠道管理优化。**

## 实现过程：

1. 锁定微博、大众点评等互联网数据源，采集百万级别消费者谈及的与该日化品牌购物相关内容；
2. 利用自然语言处理技术，对用户评论进行多维建模，包括购物环境、服务、价值等10多个一级维度和50个二级维度，实现对用户评论的量化；
3. 对沃尔玛、屈臣氏、京东等50个零售渠道进行持续监控，结果通过DashBoard和周期性分析报告呈现。

## 意义：

### 1. 有效掌握KA渠道整体情况：

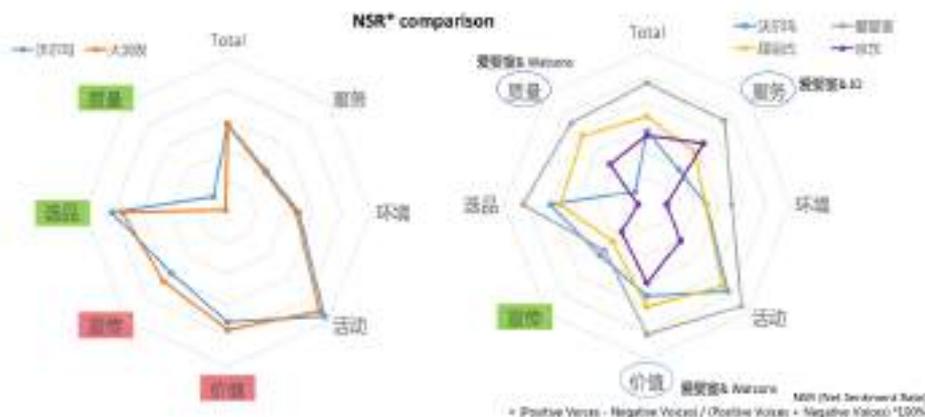
- 渠道：大卖场超市、综合商场、高端超市、化妆品店、母婴店、
- 城市：一线、二线、三线

### 2. 了解KA渠道的关键细节、优势与劣势

- 环境：货架、布局、设施、空间、位置、人气
- 服务：态度、效率、专业性

### 3. 指导渠道评级体系调整和产品促销规划

- 价值：价格、促销、会员、比价
- 选品：齐全性、独特性、进口、新品





# 大数据应用：数据驱动市场和消费者研究



## 行业

### 行业现状及趋势洞察

大数据雷达对乳制品行业生命周期及行业销量进行趋势分析，评估行业生命力与健康度；分析竞品市场占有率及与竞品区域市场比重，以寻求市场突破口或抢先产品、渠道等市场布局。

## 品牌

### 品牌健康度研究

通过大数据雷达系统从其品牌定位、品牌认知、品牌形象、品牌美誉、舆情预警、品牌资产等维度分析消费者对该乳制品品牌的认知匹配度、品牌喜好度、评价关键词正负比例、情感卷入度等，并与竞品进行对比分析，以综合评估品牌健康度。

## 产品

### 产品分析

大数据雷达系统对产品销售趋势、市场份额比重分析，判断产品所处生命周期阶段；分析消费者对产品功能、性能、包装、质量的市场关注度、认可度，为新需求发现、产品研发、定价、渠道布局、营销推广等提供策略指导。

## 消费者

### 消费者画像与需求洞察

通过大数据雷达系统分析人口属性、关系属性、内容属性；发现其潜在兴趣点和消费行为特征，以为该乳制品在营销策划、内容创作、日常运营等提供策略和个性化服务需求。

## 媒介

### 营销活动媒介策略规划

大数据雷达系统对投放媒介实行前期监测、挖掘传播KOL，并为媒介组合提供精准策略，同时根据对投放媒介进行实时跟踪，优化媒介策略和效果评估。

## 渠道

### 重客画像与监测

通过大数据雷达系统分析重复消费者线上及门店购物偏好、购物行为、购物体验、评价情感值、服务体验、购买单价等数据，按人口属性、关系属性、价值属性、经济属性等维度进行标签画像，同时实时监测个体消费者的社交行为与购买路径，以更新用户数据，丰富重客标签。

# 大数据应用：数据驱动市场和消费者研究



产品多维建模、产品优缺点分析、产品关注属性分析、竞品对比

## 产品多维度分析



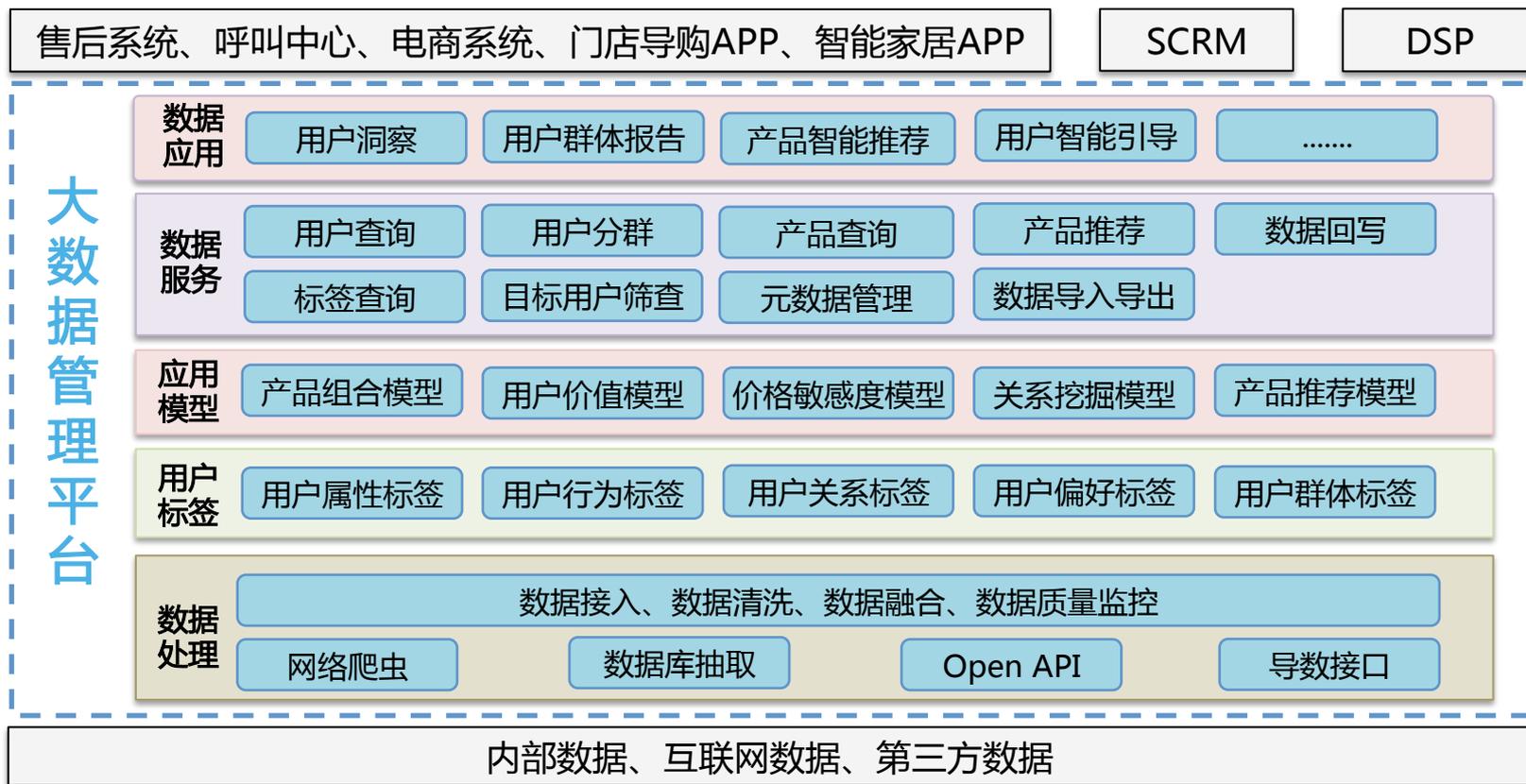
## 产品取胜点在哪？



## 用户最关注的产品属性是哪些？



# 某家电品牌大数据管理平台：应用架构



建设统一数据管理平台，实现用户全景视图并构建业务支撑应用  
将沉淀的数据资产转化为生产力

# 某家电品牌大数据管理平台：数据架构

